

Gradient-based optimization

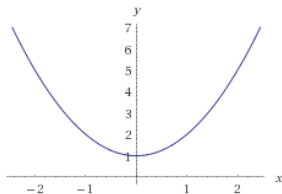
Dérivative and extrema

Theorem If f is differentiable in the vicinity of a , and has a local extremum in a , then $f'(a) = 0$.

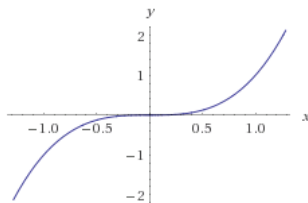
Dérivative and extrema

Theorem If f is differentiable in the vicinity of a , and has a local extremum in a , then $f'(a) = 0$.

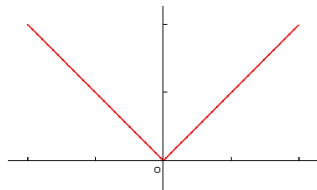
AND NOTHING MORE !!



(loc. extr. + diff.) $\implies (f'(a) = 0)$



$(f'(a) = 0) \not\Rightarrow$ loc. extr.



loc. extr. $\not\Rightarrow (f'(a) = 0)$

Necessary condition for an optimum

$$\text{Let } J : \begin{array}{ll} E \subset \mathbb{R}^n & \longrightarrow \mathbb{R} \\ \mathbf{x} = (x_1, \dots, x_n) & \longrightarrow J(x_1, \dots, x_n) \end{array}$$

- ▶ If $\hat{\mathbf{x}}$ is an internal point of E (i.e. if there exists an open set Ω such that $\hat{\mathbf{x}} \in \Omega \subset E$) and if J is differentiable in $\hat{\mathbf{x}}$, then

$$\hat{\mathbf{x}} \text{ local minimum of } J \implies \nabla J(\hat{\mathbf{x}}) = 0$$

- ▶ If E is convex, if J is convex, then

$$\hat{\mathbf{x}} \text{ local minimum of } J \iff \nabla J(\hat{\mathbf{x}}) = 0$$

Necessary condition for an optimum

$$\begin{array}{ll} \text{Let} & J : \\ & E \subset \mathbb{R}^n \longrightarrow \mathbb{R} \\ & \mathbf{x} = (x_1, \dots, x_n) \longrightarrow J(x_1, \dots, x_n) \end{array}$$

- ▶ If $\hat{\mathbf{x}}$ is an internal point of E (i.e. if there exists an open set Ω such that $\hat{\mathbf{x}} \in \Omega \subset E$) and if J is differentiable in $\hat{\mathbf{x}}$, then

$$\hat{\mathbf{x}} \text{ local minimum of } J \implies \nabla J(\hat{\mathbf{x}}) = 0$$

- ▶ If E is convex, if J is convex, then

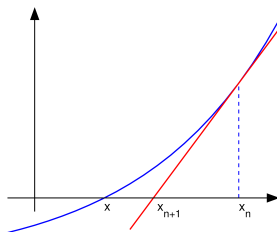
$$\hat{\mathbf{x}} \text{ local minimum of } J \iff \nabla J(\hat{\mathbf{x}}) = 0$$

Minimizing $J \implies$ finding the roots of ∇J

Descent methods

Descent methods for minimizing the cost function require the knowledge of (an estimate of) its gradient.

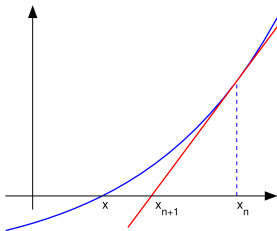
1-D: Newton iteration to find the roots of a function f : $x_{n+1} = ??$



Descent methods

Descent methods for minimizing the cost function require the knowledge of (an estimate of) its gradient.

1-D: Newton iteration to find the roots of a function f :
$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}$$

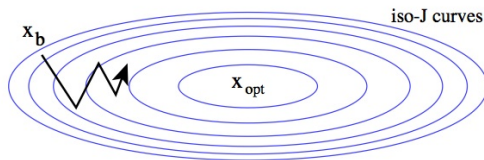


for the optimization: $f \equiv \nabla J$

Descent methods

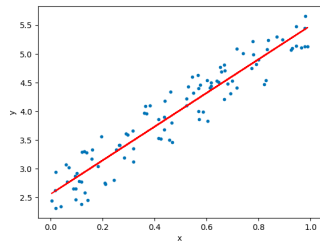
Descent methods for minimizing the cost function require the knowledge of (an estimate of) its gradient.

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{d}_k$$



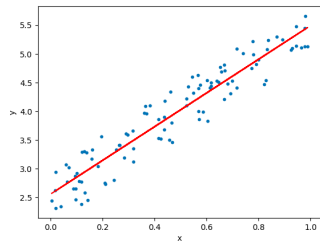
$$\text{with } \mathbf{d}_k = \begin{cases} -\nabla J(\mathbf{x}_k) & \text{gradient method} \\ -[\text{Hess}(J)(\mathbf{x}_k)]^{-1} \nabla J(\mathbf{x}_k) & \text{Newton method} \\ -\mathbf{B}_k \nabla J(\mathbf{x}_k) & \text{quasi-Newton methods (BFGS, ...)} \\ -\nabla J(\mathbf{x}_k) + \frac{\|\nabla J(\mathbf{x}_k)\|^2}{\|\nabla J(\mathbf{x}_{k-1})\|^2} \mathbf{d}_{k-1} & \text{conjugate gradient} \\ \dots & \dots \end{cases}$$

Least squares method: linear regression



Given some data $(x_i, y_i)_{i=1, \dots, p}$, what is the best approximate relationship $Y = aX + b$?

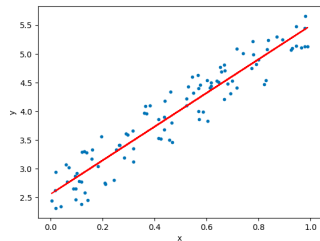
Least squares method: linear regression



Given some data $(x_i, y_i)_{i=1, \dots, p}$, what is the best approximate relationship $Y = aX + b$?

Let $y_i = ax_i + b + \varepsilon_i$ and minimize $E(a, b) = \sum_{i=1}^p \varepsilon_i^2 = \sum_{i=1}^p (y_i - ax_i - b)^2$

Least squares method: linear regression



Given some data $(x_i, y_i)_{i=1, \dots, p}$, what is the best approximate relationship $Y = aX + b$?

Let $y_i = ax_i + b + \varepsilon_i$ and minimize $E(a, b) = \sum_{i=1}^p \varepsilon_i^2 = \sum_{i=1}^p (y_i - ax_i - b)^2$

$$\begin{cases} \frac{\partial E}{\partial a}(\hat{a}, \hat{b}) = -2 \sum_{i=1}^p x_i (y_i - \hat{a} x_i - \hat{b}) = 0 \\ \frac{\partial E}{\partial b}(\hat{a}, \hat{b}) = -\sum_{i=1}^p (y_i - \hat{a} x_i - \hat{b}) = 0 \end{cases}$$

$$\text{i.e. } \begin{cases} \left(\sum_{i=1}^p x_i \right) \hat{a} + p \hat{b} = \sum_{i=1}^p y_i \\ \left(\sum_{i=1}^p x_i^2 \right) \hat{a} + \left(\sum_{i=1}^p x_i \right) \hat{b} = \sum_{i=1}^p x_i y_i \end{cases}$$

$$\text{Hence } \hat{a} = \frac{\frac{1}{p} \sum_{i=1}^p x_i y_i - \bar{x} \bar{y}}{\frac{1}{p} \sum_{i=1}^p (x_i - \bar{x})^2} \quad \text{and} \quad \hat{b} = \bar{y} - \hat{a} \bar{x}$$

Generalization: minimum of a quadratic function in finite dimension

Theorem: Generalized (or Moore-Penrose) inverse

Let \mathbf{M} a $p \times n$ matrix, with rank n , and $\mathbf{b} \in \mathbb{R}^p$.

(hence $p \geq n$)

Let $J(\mathbf{x}) = \|\mathbf{M}\mathbf{x} - \mathbf{b}\|^2 = (\mathbf{M}\mathbf{x} - \mathbf{b})^T (\mathbf{M}\mathbf{x} - \mathbf{b})$.

J is minimum for $\hat{\mathbf{x}} = \mathbf{M}^+ \mathbf{b}$, where $\mathbf{M}^+ = (\mathbf{M}^T \mathbf{M})^{-1} \mathbf{M}^T$ (generalized inverse, or Moore-Penrose inverse).

Exercise : prove it.

Generalization: minimum of a quadratic function in finite dimension

Theorem: Generalized (or Moore-Penrose) inverse

Let \mathbf{M} a $p \times n$ matrix, with rank n , and $\mathbf{b} \in \mathbb{R}^p$.

(hence $p \geq n$)

Let $J(\mathbf{x}) = \|\mathbf{M}\mathbf{x} - \mathbf{b}\|^2 = (\mathbf{M}\mathbf{x} - \mathbf{b})^T (\mathbf{M}\mathbf{x} - \mathbf{b})$.

J is minimum for $\hat{\mathbf{x}} = \mathbf{M}^+ \mathbf{b}$, where $\mathbf{M}^+ = (\mathbf{M}^T \mathbf{M})^{-1} \mathbf{M}^T$ (generalized inverse, or Moore-Penrose inverse).

Exercise : check that this is consistent with the previous results of linear regression.

Generalization: minimum of a quadratic function in finite dimension

Theorem: Generalized (or Moore-Penrose) inverse

Let \mathbf{M} a $p \times n$ matrix, with rank n , and $\mathbf{b} \in \mathbb{R}^p$.

(hence $p \geq n$)

Let $J(\mathbf{x}) = \|\mathbf{M}\mathbf{x} - \mathbf{b}\|^2 = (\mathbf{M}\mathbf{x} - \mathbf{b})^T (\mathbf{M}\mathbf{x} - \mathbf{b})$.

J is minimum for $\hat{\mathbf{x}} = \mathbf{M}^+ \mathbf{b}$, where $\mathbf{M}^+ = (\mathbf{M}^T \mathbf{M})^{-1} \mathbf{M}^T$ (generalized inverse, or Moore-Penrose inverse).

Corollary: with a generalized norm

Let \mathbf{N} a $p \times p$ symmetric definite positive matrix.

Let $J_1(\mathbf{x}) = \|\mathbf{M}\mathbf{x} - \mathbf{b}\|_{\mathbf{N}}^2 = (\mathbf{M}\mathbf{x} - \mathbf{b})^T \mathbf{N} (\mathbf{M}\mathbf{x} - \mathbf{b})$.

J_1 is minimum for $\hat{\mathbf{x}} = (\mathbf{M}^T \mathbf{N} \mathbf{M})^{-1} \mathbf{M}^T \mathbf{N} \mathbf{b}$.

Just for scientific culture: optimization in infinite dimension, using the adjoint

$$\begin{cases} -u''(x) + c(x) u'(x) = f(x) & x \in]0, 1[, f \in L^2(]0, 1[) \\ u(0) = u(1) = 0 \end{cases}$$

- ▶ $c(x)$ is **unknown**
- ▶ $u^{\text{obs}}(x)$ is an **observation** of $u(x)$

Problem How can we estimate the function $c(x)$?

Just for scientific culture: optimization in infinite dimension, using the adjoint

$$\begin{cases} -u''(x) + c(x) u'(x) = f(x) & x \in]0, 1[, f \in L^2(]0, 1[) \\ u(0) = u(1) = 0 \end{cases}$$

- ▶ $c(x)$ is **unknown**
- ▶ $u^{\text{obs}}(x)$ is an **observation** of $u(x)$

Problem How can we estimate the function $c(x)$?

Idea minimize the **cost function**: $J(c) = \frac{1}{2} \int_0^1 (u(x) - u^{\text{obs}}(x))^2 dx$

→ need to compute $\nabla J(c)$

Example of the adjoint method

$$\begin{cases} -u''(x) + c(x) u'(x) = f(x) & x \in]0, 1[, f \in L^2(]0, 1[) \\ u(0) = u(1) = 0 \end{cases}$$

Find $c(x)$ that minimizes $J(c) = \frac{1}{2} \int_0^1 \left(u(x) - u^{\text{obs}}(x) \right)^2 dx$

Example of the adjoint method

$$\begin{cases} -u''(x) + c(x) u'(x) = f(x) & x \in]0, 1[, f \in L^2(]0, 1[) \\ u(0) = u(1) = 0 \end{cases}$$

Find $c(x)$ that minimizes $J(c) = \frac{1}{2} \int_0^1 \left(u(x) - u^{\text{obs}}(x) \right)^2 dx$

$\nabla J \rightarrow$ Gâteaux-derivative: $\hat{J}[c](\delta c) = \langle \nabla J(c), \delta c \rangle$

$$\hat{J}[c](\delta c) = \int_0^1 \hat{u}(x) \left(u(x) - u^{\text{obs}}(x) \right) dx \quad \text{with } \hat{u} = \lim_{\alpha \rightarrow 0} \frac{u_{c+\alpha \delta c} - u_c}{\alpha}$$

What is the equation satisfied by \hat{u} ?

$$\begin{cases} -\hat{u}''(x) + c(x) \hat{u}'(x) = -\delta c(x) u'(x) & x \in]0, 1[\\ \hat{u}(0) = \hat{u}(1) = 0 \end{cases} \quad \begin{array}{l} \text{tangent} \\ \text{linear model} \end{array}$$

Example of the adjoint method

$$\begin{cases} -u''(x) + c(x) u'(x) = f(x) & x \in]0, 1[, f \in L^2(]0, 1[) \\ u(0) = u(1) = 0 \end{cases}$$

Find $c(x)$ that minimizes $J(c) = \frac{1}{2} \int_0^1 (u(x) - u^{\text{obs}}(x))^2 dx$

$\nabla J \rightarrow$ Gâteaux-derivative: $\hat{J}[c](\delta c) = \langle \nabla J(c), \delta c \rangle$

$$\hat{J}[c](\delta c) = \int_0^1 \hat{u}(x) (u(x) - u^{\text{obs}}(x)) dx \quad \text{with } \hat{u} = \lim_{\alpha \rightarrow 0} \frac{u_{c+\alpha \delta c} - u_c}{\alpha}$$

What is the equation satisfied by \hat{u} ?

$$\begin{cases} -\hat{u}''(x) + c(x) \hat{u}'(x) = -\delta c(x) u'(x) & x \in]0, 1[\\ \hat{u}(0) = \hat{u}(1) = 0 \end{cases} \quad \begin{array}{l} \text{tangent} \\ \text{linear model} \end{array}$$

Going back to \hat{J} scalar product of the TLM with a variable p : $-\int_0^1 \hat{u}'' p + \int_0^1 c \hat{u}' p = -\int_0^1 \delta c u' p$

Integration by parts: $\int_0^1 \hat{u} (-p'' - (c p)') = \hat{u}'(1)p(1) - \hat{u}'(0)p(0) - \int_0^1 \delta c u' p$

Example of the adjoint method

$$\begin{cases} -u''(x) + c(x) u'(x) = f(x) & x \in]0, 1[, f \in L^2(]0, 1[) \\ u(0) = u(1) = 0 \end{cases}$$

Find $c(x)$ that minimizes $J(c) = \frac{1}{2} \int_0^1 (u(x) - u^{\text{obs}}(x))^2 dx$

$\nabla J \rightarrow$ Gâteaux-derivative: $\hat{J}[c](\delta c) = \langle \nabla J(c), \delta c \rangle$

$$\hat{J}[c](\delta c) = \int_0^1 \hat{u}(x) (u(x) - u^{\text{obs}}(x)) dx \quad \text{with } \hat{u} = \lim_{\alpha \rightarrow 0} \frac{u_{c+\alpha \delta c} - u_c}{\alpha}$$

What is the equation satisfied by \hat{u} ?

$$\begin{cases} -\hat{u}''(x) + c(x) \hat{u}'(x) = -\delta c(x) u'(x) & x \in]0, 1[\\ \hat{u}(0) = \hat{u}(1) = 0 \end{cases} \quad \begin{matrix} \text{tangent} \\ \text{linear model} \end{matrix}$$

Going back to \hat{J} scalar product of the TLM with a variable p : $-\int_0^1 \hat{u}'' p + \int_0^1 c \hat{u}' p = -\int_0^1 \delta c u' p$

Integration by parts: $\int_0^1 \hat{u} (-p'' - (c p)') = \hat{u}'(1)p(1) - \hat{u}'(0)p(0) - \int_0^1 \delta c u' p$

$$\begin{cases} -p''(x) - (c(x) p(x))' = u(x) - u^{\text{obs}}(x) & x \in]0, 1[\\ p(0) = p(1) = 0 \end{cases} \quad \begin{matrix} \text{adjoint} \\ \text{model} \end{matrix}$$

Conclusion: $\nabla J(c(x)) = -u'(x) p(x)$